

## DOCUMENT RESUME

ED 448 174

TM 032 135

AUTHOR Tanguma, Jesus  
TITLE A Review of the Literature on Missing Data.  
PUB DATE 2000-11-16  
NOTE 24p.; Paper presented at the Annual Meeting of the Mid-South Educational Research Association (28th, Bowling Green, KY, November 15-17, 2000).  
PUB TYPE Information Analyses (070) -- Numerical/Quantitative Data (110) -- Speeches/Meeting Papers (150)  
EDRS PRICE MF01/PC01 Plus Postage.  
DESCRIPTORS Correlation; \*Data Analysis; Literature Reviews; \*Regression (Statistics); \*Research Methodology  
IDENTIFIERS \*Imputation; \*Missing Data

## ABSTRACT

This paper reviews the literature on methods for dealing with missing data, discusses four commonly used methods, and illustrates these approaches with a small hypothetical data set. Most studies contain some missing data, and the reasons data are missing are many and varied. Four commonly used methods have been identified in the literature: (1) listwise deletion; (2) pairwise deletion; (3) mean imputation; and (4) regression imputation. Listwise deletion, which is the default in some statistical packages (e.g., the Statistical Package for the Social Sciences and the Statistical Analysis System), is the most commonly used method, also by default. However, because listwise deletion eliminates all cases for a participant missing data on any predictor or criterion variable, it is not the most effective method. Pairwise deletion uses those observations that have no missing values to compute the correlations. Thus, it preserves information that would have been lost when using listwise deletion. However, since different sample sizes go into the computing of the correlations, the resulting correlation matrix may not be positive definite (a mathematical condition required to invert the correlation matrix). In mean imputation, the mean for a particular variable, computed from available cases, is substituted in place of missing data values on the remaining cases. This allows the researcher to use the rest of the participant's data. When using a regression-based procedure to estimate the missing values, the estimation takes into account the relationships among the variables. Thus, substitution by regression is more statistically efficient. (Contains 1 figure, 7 tables, and 15 references.) (Author/SLD)

Reproductions supplied by EDRS are the best that can be made  
from the original document.

Running head: MISSING DATA

ED 448 174

## A Review of the Literature on Missing Data

Jesus Tanguma

University of Houston Clear Lake

U.S. DEPARTMENT OF EDUCATION  
EDUCATIONAL RESOURCES INFORMATION  
CENTER (ERIC)

- ☒ This document has been reproduced as received from the person or organization originating it.
- ☐ Minor changes have been made to improve reproduction quality.
- Points of view or opinions stated in this document do not necessarily represent official OERI position or policy.

J. Tanguma

1

Paper presented at the annual meeting of the Mid-South Educational Research

Association, Bowling Green, KY, November 16, 2000.

BEST COPY AVAILABLE

TM032135

## Abstract

Most studies contain some missing data. The reasons for the missing data are many and varied. Respondents did not provide complete information. Observers failed to record all pertinent information. Participants did not participate throughout the duration of the study. Data was not properly coded/transferred.

Four commonly used methods (listwise deletion, pairwise deletion, mean imputation, and regression imputation) for dealing with missing data are illustrated by means of a hypothetical example.

Listwise deletion, being the default in some statistical packages (e.g., SPSS and SAS), is the one most commonly used method, also by default. However, because listwise deletion eliminates all cases for a participant missing data on any predictor or criterion variable, it is not the most effective method.

Pairwise deletion uses those observations that have no missing values to compute the correlations. Thus, it preserves information that would have been lost when using listwise deletion. However, since different sample sizes go into the computing of the correlations, the resulting correlation matrix may not be positive definite (a mathematical condition required to invert the correlation matrix).

In mean imputation, the mean for a particular variable, computed from available cases, is substituted in place of missing data values on the remaining cases. This allows the researcher to use the rest of the participant's data.

When using a regression-based procedure to estimate the missing values, the estimation takes into account the relationships among the variables. Thus, substitution by regression is more statistically efficient.

### A Review of the Literature on Missing Data

Most research studies (e.g., survey studies and field experiments) contain some missing data. However, most standard statistical methods have been designed to analyze data sets with no missing data. Consequently, the researcher has two options (a) to delete those cases which have missing data, or (b) to fill-in the missing values with estimated values (Anderson, Basilevsky, & Hum, 1983). Thus, a data set is created containing no missing values (empty cells). Typically, the data set is presented in a rectangular table where rows indicate cases, observations, or subjects, and columns indicate variables measured on each unit (Little & Rubin, 1987).

The reasons for the missing data are many and varied. Respondents did not provide complete information. Observers failed to record all pertinent information. Participants did not participate throughout the duration of the study. Data was not properly coded/transferred. Data/instrument was lost. The fact of the matter is that, as so eloquently stated by Cohen and Cohen (1983), "if there are any ways in which data can be missing, they will be" (p. 275).

There exist a number of statistical techniques (e.g., listwise deletion, pairwise deletion, mean imputation, regression imputation, hot-deck imputation, expectation maximization, and so on) for researchers to use when faced with missing data. The most obvious option is to simply drop any case that may have any missing data. For example, when a participant does not answer any of the items in the survey, that participant should not be included in the data analysis. However, this would restrict the extent to which the sample is a representative of the original population. Thus, limiting the generalizability of

the study. On the other hand, when the participant partially answers the survey, the question is whether or not to include the subject in the data analysis. If the subject's data enters into the analysis, how should the missing data be handled? Before deciding on this, it might be instructive to see if the data is missing on the dependent or the independent variables.

Cohen and Cohen (1983) have suggested that when the missing data is on the dependent variable, the subject may be dropped from the analysis. However, if the missing data is among the independent variables, it might be instructive to determine what proportion of the data is missing. According to Orme and Reis (1991) "if a large proportion of data is missing, the validity of the study can be so compromised that it would be best to redesign the study and conduct it again" (p. 62). On the other hand, if only a small to moderate proportion of the data is missing for one or several independent variables, the different techniques to handle missing data may lead to different results. Thus, causing confusion to the applied researcher. This may be the case, for example, when the researcher allows the computer package to use the default options. However, since some computer packages (e.g., SPSS and SAS) have listwise and pairwise deletion (depending on the applications) as their defaults, the uniformed researcher will be using listwise or pairwise deletion methods, also by default. However, as emphasized by the APA Task Force on Statistical Inference on their recently released report:

Special issues arise in modeling when we have missing data. The two popular methods for dealing with missing data that are in basic statistics packages—listwise and pairwise deletion of missing values—are among the worst methods available for practical

applications. (Wilkinson & APA Task Force on Statistical Inference, 1999, p. 598)

The purpose of the present paper is to discuss and illustrate four commonly used methods (listwise deletion, pairwise deletion, mean substitution, and regression imputation) for dealing with missing data. To make the discussion and illustration more concrete, a small hypothetical data set will be used. The interested reader may recompute the results using the heuristic data set and thus obtain a better understanding of the methods and procedures presented.

#### Listwise Deletion

Listwise deletion drops any case on which any variable is missing any data. In doing so, any subsequent calculations/computations (e., correlation matrix, regression beta weights) are performed using a sample size somewhat smaller than the one intended. For example, after randomly deleting six entries from Table 1, the correlation matrix and regression beta weights are computed using a sample size of  $n = 14$  instead of the original  $n = 20$ . In other words, there is a 4.5% loss of data, see Table 2. Thus, listwise deletion sacrifices a large amount of data (Malhorta, 1987; Stumpf, 1978). The large loss of data will reduce the statistical power (Cohen & Cohen, 1983; Gilley & Leone, 1991) and may reduce the precision of the parameters being estimated (Cohen & Cohen, 1983; Donner, 1982; Little & Rubin, 1987). Additionally, when the data are missing at random, "type II error rates may be artificially inflated" (Raymond, 1986, p. 399). Thus, listwise deletion is not a generally adequate method for handling the missing data problem (Cohen & Cohen, 1983). However, unless specifically instructed by the researcher, SPSS and SAS will use the listwise deletion method for handling missing data, their default option.

---

Insert Tables 1 and 2 about here

---

The means and standard deviations for the original data set and those computed after using the listwise deletion method are presented in Table 3. Notice that, since X1 had no missing values (see Table 2), its mean and standard deviation remained constant. However, all other variables had different means and standard deviations as a result of deleting some cases. A pictorial representation of the different mean values is shown in Figure 1.

---

Insert Table3 about here

---

Just as the means and standard deviations of the predictor variables changed after deleting some case values, so did the unstandardized regression coefficients, see Table 4. For example, the unstandardized regression coefficient for X2 when using the original data set is 0.538. However, after deleting some cases the unstandardized regression coefficient for X2 is now 0.708. Thus, using listwise deletion to predict some outcome variable when some of the predictors contain missing data does affect the unstandardized regression coefficients.

---

Insert Table 4 about here

---

#### Pairwise Deletion

Pairwise deletion computes means, variances and standard deviations from available cases. The correlation coefficients are computed from all cases with values on the (two) variables involved. As shown in Table 5, the correlation coefficients obtained

for the original data set differ from those obtained applying the pairwise deletion method to the data set with missing data. Another interesting point from Table 5 is that the sample sizes on which the different pairwise correlations are computed vary. Thus, making it unclear as to what sample size to use for the computation of standard errors and tests of statistical significance (Orme & Reis, 1991). In addition, the different sample sizes on which the pairwise correlations are computed make the population to which one can generalize somewhat unclear. Other problems associated with the use of pairwise deletion are that the correlations being estimated may lie outside the acceptable range (-1, 1) and that the  $R^2$  may be less than zero or larger than one (Cohen & Cohen, 1983; Raymond, 1987; Little & Rubin, 1987). Additionally, as pointed out by Kim and Curry (1977), "the matrix generated by pairwise deletion may not be consistent (not positive definite), especially when the missing data pattern is not random or when the total sample size is small" (p. 222). Positive definite is a mathematical condition required to invert the correlation matrix. If the correlation matrix can not be inverted, this can have serious negative effects on maximum likelihood-based programs such as AMOS, LISREL and PROC CALIS in SAS (Roth, 1994).

#### Imputation Methods

The previous two (listwise and pairwise) methods of handling missing data make use of the data that are available only. However, in some instances it might be prudent to fill-in (impute) the missing cases. By imputing the missing values, the researcher is then able to use standard statistical techniques that require complete data sets. Additionally, the recovery of sample size and statistical power is a motivational factor in imputing values (Raymond, 1987). Although a variety of methods for estimating (imputing)



missing values have been proposed, only two techniques will be presented in the following sections.

Imputation of missing values by sensible estimates, although widely used, has some pitfalls (Little & Rubin, 1987). According to Dempster and Rubin (1983):

The idea of imputation is both seductive and dangerous. It is seductive because it can lull the user into the pleasure state of believing that the data are complete after all, and it is dangerous because it lumps together situations where the problem is sufficiently minor that it can be legitimately handled in this way and situations where standard estimators applied to the real imputed data have substantial biases.

#### Mean Imputation

According to Raymond (1986), "the most widely used estimation technique is probably the mean substitution method" (p. 403). By filling-in the missing cases, the researcher restores the sample size to its original size. However, because the means are replacing the missing values, variances and covariances will be downwardly biased (Little & Rubin, 1987). Recall that a formula for computing the variance for a sample is

$s^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$ . Thus, when some of the  $X_i$ 's (raw scores) have been replaced by the

mean ( $\bar{X}$ ) of the distribution, the sum of squares does not change. In other words, only zeros are being added to the sum of squares obtained when there were missing values. Yet, the sample size ( $n$ ) has increased. Consequently, the variance will be decreased. For example, the variance for variable X7 after mean imputation is 6.532. However, the

variance for the same variable using the original data set is 7.082. Again, this is because the numerator of the variance formula did not change but the denominator did increase.

Another problem with mean imputation is that the correlation coefficients are attenuated. A formula for computing the correlation coefficient between two variables is

$$r_{xy} = \frac{\sum z_x z_y}{n-1}. \text{ But given that } z_x = \frac{X_i - \bar{X}}{s_x}, \text{ it follows that when } X_i = \bar{X}, z_x \text{ is not}$$

contributing to the summation. Therefore, when the imputation has been done by the means, it follows that the numerator of the formula for computing the correlation coefficients does not change yet the sample size does increase. Thus, the correlation coefficients under mean imputation will be downwardly biased (Raymond, 1986). For example, the correlation coefficient between X2 and X3 under the mean imputation is 0.097. On the other hand, the correlation coefficient between X2 and X3 computed from the original data set is 0.113.

Just as variances and covariances are attenuated when imputing by the means, the confidence intervals may not be as precise as expected. As Little and Rubin (1990) have pointed out

95% confidence intervals for parameters computed from the filled-in data may in fact cover the true parameter value only 80% to 90% of the time, and tests with nominal significance level of 5% may have a true significance level of 10% or 20%. (p. 294)

### Regression Imputation

The second imputation technique to be discussed in this section is regression imputation. Although other regression imputation techniques exist (e.g., stepwise or iterative regression), only the simplest case (single iteration) will be illustrated here. The

imputed data will preserve deviations from the mean as well as the shape of the distribution (Little, 1988). Thus, according to Roth (1994), the imputed data "will not attenuate correlations as much as mean substitution" (p. 542).

Regression imputation is done in several steps. To better illustrate the procedure, the data set in Table 2 will be used. Notice that Table 2 has some empty cells. These are the cells to be imputed by regression. For example, to compute the missing value for variable  $X_i$ , all other variables (i.e.,  $X_j$ , where  $j = 1, \dots, 7$  but  $i \neq j$ ) are regressed (ironically, this is often done using listwise deletion) on the variable of interest ( $X_i$ ). Next, the regression weight(s) are applied to the known scores for  $X_i$  to calculate the value for the empty cell. Symbolically,

$$X_i = \hat{y} = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4 + b_5 X_5 + b_6 X_6 + b_7 X_7.$$

Table 6 presents the different regression weights used to impute the missing values. For example, in imputing the missing value for  $X_2$ , the following compute statement would be used:

$$X_2 = 1.003 - 0.693X_1 + 0.266X_3 + 0.630X_4 + 0.160X_5 + 0.112X_6 + 0.174X_7.$$

Thus,  $X_2 = 6.24$ . The rest of the replace/imputed values are presented in Table 7.

---

Insert Tables 6 and 7 about here

---

When regression imputation is used to fill-in missing values on the dependent variable, those with missing values on the dependent variables will be perfectly predicted. Thus, inflating the predictive power of the model. On the other hand, if regression imputation is used to fill-in missing values on the independent variables, the imputed

values will be perfectly correlated with the other variables in the model. Thus, increasing multicollinearity among the independent variables.

### Conclusion

Missing data are a common problem in most research studies. Yet no commonly agreed upon solution exists. Consequently, researchers have developed a wide variety of techniques for handling missing data. However, no single technique is without pitfalls. Thus, researchers facing a missing data problem should thoroughly investigate the sources of the missing data as well as the options for handling missing data.

This paper has presented four techniques for handling missing data. When using listwise deletion, there is a large loss of subjects/cases. This loss of data, will reduce the statistical power, may reduce the precision of the parameters being investigated, and may inflate the Type II error rates. Thus, listwise deletion is not a generally adequate method for handling the missing data problem (Cohen & Cohen, 1986). However, since this is the default in most statistical packages (e.g., SPSS, and SAS), listwise deletion will probably continue to be used, also by default.

Using pairwise deletion would save some of the data that would be lost if listwise deletion would be used. However, because the sample sizes on which the different pairwise correlations are computed vary, it is unclear what sample size to use in the computation of standard errors and tests of statistical significance (Orme & Reis, 1991). Thus, posing a potential threat to statistical conclusion validity. Additionally, the matrix generated by pairwise deletion may not be positive definite.

Mean imputation will restore the sample size to its original size. However, because the means are replacing the missing values, variances and covariances will be downwardly biased (Little & Rubin, 1987). Additionally, the confidence intervals may not be as precise as expected (Little & Rubin, 1990).

Although regression imputation "will not attenuate correlations as much as mean substitution" (Roth, 1994, p. 542), the method is not without pitfalls. Imputing missing values on dependent variables by regression will inflate the predictive power of the model. Imputing missing values on the independent variables will increase multicollinearity.

The intent of this paper has been to alert applied researchers as to what effects do the different techniques for dealing with missing data have on parameters, variances, correlations and confidence intervals. By working through the examples, the applied researcher might realize that perhaps it is best not to use the defaults on some of the statistical packages (e.g., SPSS and SAS). Instead, the applied researcher should thoroughly investigate the available options before deciding on a specific technique for handling missing data.

## References

- Anderson, A. B., Basilevsky, A. & Hum, D. P. J. (1983). Missing data: A review of the literature. In P. H. Rossi, J. D. Wright, & A. B. Anderson (Eds.), Handbook of survey research (pp. 415-494). San Diego: Academic Press.
- Cohen, J. & Cohen, P. (1983). Missing data. In J. Cohen & P. Cohen, Applied multiple regression: Correlation analysis for the behavioral sciences (pp. 275-300).
- Dempster, A. P. & Rubi, D. B. (1983). Overview, in Incomplete Data in sample Surveys, Vol 2: Theory and Annotated Bibliography (W. G. Madow, I. Olkin, & D. B. Rubin, Eds.). New York: Academic Press, 3-10.
- Donner, A. (1982). The relative effectiveness of procedures commonly used in multiple regression analysis for dealing with missing data. The American Statistician, 36, 378-381.
- Gilley, O. W. & Leone, R. P. (1991). A two-stage imputation procedure for item nonresponse in surveys. Journal of Business Research, 22, 281-291.
- Kim, J. O. & Curry, J. (1977). The treatment of missing data in multivariate analysis. Sociological Methods & Research, 6, 215-241.
- Little, R. J. A. (1990). Missing data adjustments in large surveys. Journal of Business & Economics Statistics, 6, 1-15
- Little, R. J. A. & Rubin, D. R. (1987). Statistical analysis with missing data. New York: Wiley.
- Little, R. J. A. & Rubin, D. R. (1990). The analysis of social science data with missing values. Sociological Methods & Research, 18, 292-326.

Malhorta, N. K. (1987). Analyzing marketing research data with incomplete information on the dependent variable. Journal of Marketing Research, 24, 74-84.

Orme, J. G. & Reis, J. (1991). Multiple regression with missing data. Journal of Social Service Research, 15, 61-91.

Raymond, M. R. (1986). Missing data in evaluation research. Evaluation & the Health Profession, 9, 395-420.

Roth, P. L. (1994). Missing data: A conceptual review for applied psychologists. Personnel Psychology, 47, 537-560.

Stumpf, S. A. (1978). A note on handling missing data. Journal of Management, 4, 65-73.

Wilkinson, L. & The APA Task Force on Statistical Inference. (1999). Statistical methods in psychology journals: Guidelines and explanations. American Psychologist, 54, 594-604.

Figure 1. Means across methods

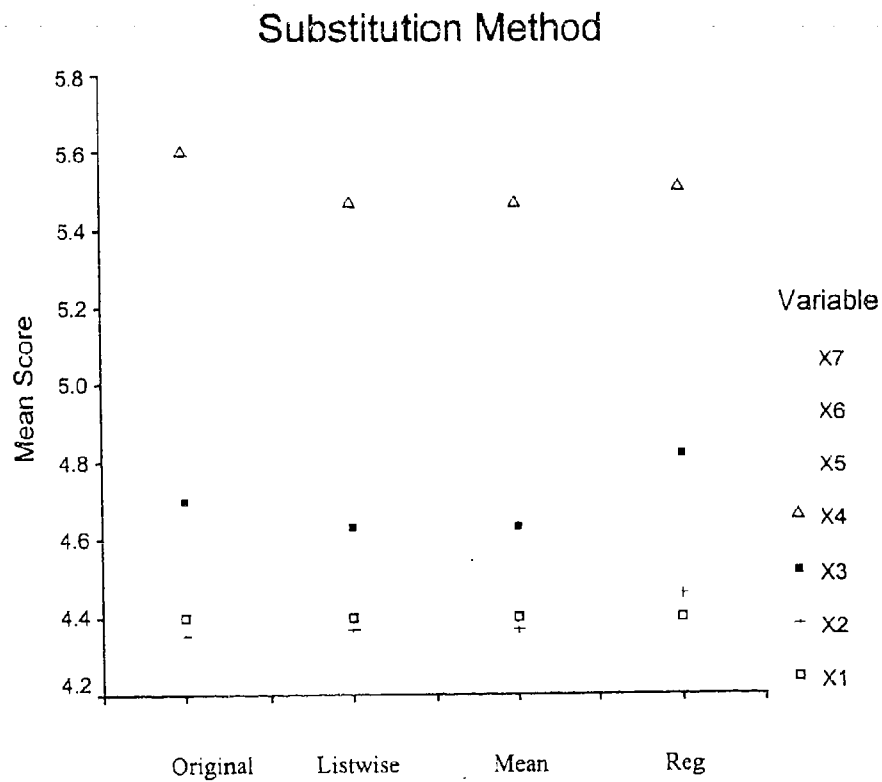




Table 1. Original data set

Y	X1	X2	X3	X4	X5	X6	X7
2	7	5	6	8	2	7	3
4	1	4	3	4	8	1	7
6	2	2	6	1	7	3	3
3	2	3	5	8	2	7	8
2	6	4	3	8	6	5	9
7	7	1	5	4	5	4	8
4	2	7	2	8	4	3	1
5	2	8	3	8	4	7	8
1	7	2	4	7	5	4	6
2	7	2	7	8	3	6	2
2	6	3	2	3	8	9	7
1	2	2	2	3	9	8	3
6	4	8	8	6	3	8	7
2	4	2	4	1	4	2	1
7	1	6	8	4	1	5	5
8	8	4	6	7	5	2	4
6	2	8	3	5	5	8	1
3	3	5	7	5	4	5	3
7	8	3	4	6	8	1	4
1	7	8	6	8	8	9	7

Table 2. Data set with missing values

Y	X1	X2	X3	X4	X5	X6	X7
2	7	5		8	2	7	3
4	1		3	4	8	1	7
6	2	2	6	1		3	3
3	2	3	5	8	2		8
2	6	4	3		6	5	9
7	7	1	5	4	5	4	
4	2	7	2	8	4	3	1
5	2	8	3	8	4	7	8
1	7	2	4	7	5	4	6
2	7	2	7	8	3	6	2
2	6	3	2	3	8	9	7
1	2	2	2	3	9	8	3
6	4	8	8	6	3	8	7
2	4	2	4	1	4	2	1
7	1	6	8	4	1	5	5
8	8	4	6	7	5	2	4
6	2	8	3	5	5	8	1
3	3	5	7	5	4	5	3
7	8	3	4	6	8	1	4
1	7	8	6	8	8	9	7

Table 3. Means and standard deviations for various data sets

Variable	Original		Listwise		Mean		Regression	
	Data set		Deletion		Substitution		Substitution	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
X1	4.40	2.56	4.40	2.56	4.40	2.56	4.40	2.56
X2	4.35	2.39	4.37	2.45	4.37	2.39	4.46	2.43
X3	4.70	1.98	4.63	2.01	4.63	1.95	4.82	2.13
X4	5.60	2.39	5.47	2.39	5.47	2.33	5.51	2.33
X5	5.05	2.33	4.95	2.34	4.95	2.29	4.85	2.32
X6	5.20	2.63	5.11	2.66	5.11	2.59	5.16	2.60
X7	4.85	2.66	4.68	2.63	4.68	2.56	4.63	2.57

Table 4. Unstandardized regression coefficients for various data sets

	const	X1	X2	X3	X4	X5	X6	X7
Original	4.989	.138	.538	.197	-.409	-.172	-.504	.179
Listwise	3.910	.117	.708	.128	-.270	-.064	-.570	.103
Mean	4.547	.146	.530	.266	-.337	-.115	-.482	.019
Regression	6.077	.232	.510	.033	-.370	-.306	-.421	.024

Table 5. Correlation coefficients for various data sets

Variable		X1	X2	X3	X4	X5	X6	X7
X1	Original	1.000	-.290	.150	.336	.138	-.052	.125
	Pairwise	1.000	-.317	.117	.314	.189	-.017	.063
	n		19	19	19	19	19	19
	Mean	1.000	-.301	.114	.310	.185	-.016	.061
	Reg	1.000	-.351	.199	.321	.223	-.037	.040
X2	Original	-.290	1.000	.113	.412	-.230	.399	.034
	Pairwise	-.317	1.000	.100	.433	-.188	.447	.154
	n	19	19	18	18	18	18	18
	Mean	-.301	1.000	.097	.428	-.178	.413	.138
	Reg	-.351	1.000	.077	.391	-.073	.328	.201
X3	Original	.150	.113	1.000	.085	-.546	.043	-.019
	Pairwise n	.117	.100	1.000	.106	-.580	.011	-.006
		19	18	19	18	18	18	18
	Mean	.114	.097	1.000	.098	-.547	.011	-.006
	Reg	.199	.077	1.000	.174	-.636	.079	-.069
X4	Original	.336	.412	.085	1.000	-.327	.223	.255
	Pairwise n	.314	.433	.106	1.000	-.308	.202	.250
		19	18	18	19	18	18	18
	Mean	.310	.428	.098	1.000	-.276	.195	.223
	Reg	.321	.391	.174	1.000	-.175	.216	.266

Table 5 continued

X5	Original	.138	-.230	-.546	-.327	1.000	-.079	.112
	Pairwise	.189	-.188	-.580	-.308	1.000	.005	.153
	n	19	18	18	18	19	18	18
	Mean	.185	-.178	-.547	-.276	1.000	.008	.151
	Reg	.223	-.073	-.636	-.175	1.000	.017	.176
X6	Original	-.052	.399	.043	.223	-.079	1.000	.215
	Pairwise n	-.017	.447	.011	.202	.005	1.000	.222
		19	18	18	18	18	19	18
	Mean	-.016	.413	.011	.195	.008	1.000	.209
	Reg	-.037	.328	.079	.216	.017	1.000	.245
X7	Original	.125	.034	-.019	.255	.112	.215	1.000
	Pairwise	.063	.154	-.006	.250	.153	.222	1.000
	n	19	18	18	18	18	18	19
	Mean	.061	.138	-.006	.223	.151	.209	1.000
	Reg	.040	.201	-.069	.266	.176	.245	1.000

Table 6. Regression weights used to compute missing values

Predict X2							
	a	X1	X3	X4	X5	X6	X7
Unst B coeff	1.003	-.693	.266	.630	.160	.112	.174
Predict X3							
	a	X1	X2	X4	X5	X6	X7
Unst B coeff	5.378	.609	.215	-.250	-.817	.171	.004
Predict X4							
	a	X1	X2	X3	X5	X6	X7
Unst B coeff	2.947	.701	.597	-.292	-.400	.011	-.003
Predict X5							
	a	X1	X2	X3	X4	X6	X7
Unst B coeff	4.923	.625	.105	-.660	-.276	.276	-.003
Predict X6							
	a	X1	X2	X3	X4	X5	X7
Unst B coeff	.336	-.560	.179	.337	.019	.671	.407
Predict X7							
	a	X1	X2	X3	X4	X5	X6
Unst B coeff	-1.438	.423	.268	.074	-.005	-.008	.394

Table 7. Data set with missing values replaced

X1m	X2m	X3m	X4m	X5m	X6m	X7m	X1re	X2re	X3re	X4re	X5re	X6re	X7re
7.00	5.00	4.63	8.00	2.00	7.00	3.00	7.00	5.00	8.39	8.00	2.00	7.00	3.00
1.00	4.37	3.00	4.00	8.00	1.00	7.00	1.00	6.24	3.00	4.00	8.00	1.00	7.00
2.00	2.00	6.00	1.00	4.95	3.00	3.00	2.00	2.00	6.00	1.00	2.97	3.00	3.00
2.00	3.00	5.00	8.00	2.00	5.11	8.00	2.00	3.00	5.00	8.00	2.00	6.19	8.00
6.00	4.00	3.00	5.47	6.00	5.00	9.00	6.00	4.00	3.00	6.30	6.00	5.00	9.00
7.00	1.00	5.00	4.00	5.00	4.00	4.67	7.00	1.00	5.00	4.00	5.00	4.00	3.68
2.00	7.00	2.00	8.00	4.00	3.00	1.00	2.00	7.00	2.00	8.00	4.00	3.00	1.00
2.00	8.00	3.00	8.00	4.00	7.00	8.00	2.00	8.00	3.00	8.00	4.00	7.00	8.00
7.00	2.00	4.00	7.00	5.00	4.00	6.00	7.00	2.00	4.00	7.00	5.00	4.00	6.00
7.00	2.00	7.00	8.00	3.00	6.00	2.00	7.00	2.00	7.00	8.00	3.00	6.00	2.00
6.00	3.00	2.00	3.00	8.00	9.00	7.00	6.00	3.00	2.00	3.00	8.00	9.00	7.00
2.00	2.00	2.00	3.00	9.00	8.00	3.00	2.00	2.00	2.00	3.00	9.00	8.00	3.00
4.00	8.00	8.00	6.00	3.00	8.00	7.00	4.00	8.00	8.00	6.00	3.00	8.00	7.00
4.00	2.00	4.00	1.00	4.00	2.00	1.00	4.00	2.00	4.00	1.00	4.00	2.00	1.00
1.00	6.00	8.00	4.00	1.00	5.00	5.00	1.00	6.00	8.00	4.00	1.00	5.00	5.00
8.00	4.00	6.00	7.00	5.00	2.00	4.00	8.00	4.00	6.00	7.00	5.00	2.00	4.00
2.00	8.00	3.00	5.00	5.00	8.00	1.00	2.00	8.00	3.00	5.00	5.00	8.00	1.00
3.00	5.00	7.00	5.00	4.00	5.00	3.00	3.00	5.00	7.00	5.00	4.00	5.00	3.00
8.00	3.00	4.00	6.00	8.00	1.00	4.00	8.00	3.00	4.00	6.00	8.00	1.00	4.00
7.00	8.00	6.00	8.00	8.00	9.00	7.00	7.00	8.00	6.00	8.00	8.00	9.00	7.00